

Stochastic Gradient Descent for Numerical MAX-CSP

Derek Paulsen

1 Abstract

2 Introduction

3 Background

3.0.1 Constraint satisfaction problems

Constraint satisfaction problems (CSPs) encompasses many different problems. Formally a constraint satisfaction problem (CSP) is defined as a triple $\langle X, D, C \rangle$, where

$X = \{X_1, \dots, X_n\}$ the set of variables

$D = \{D_1, \dots, D_n\}$ the domains of each respective variable

$C = \{C_1, \dots, C_m\}$ the set of constraints

A classic example of CSP, is the 3-SAT problem. In this case the boolean variables that appear in at least one clause would be X . The domain for each variable would be $D_i = \{true, false\}$ and the constraints would all be three variable disjunctive clauses. The 3-SAT problem of course known to be NP-Complete. In this paper we focus on a subset of CSPs called Numerical MAX-CSP. In typical a CSP all constraints must be satisfied, however in MAX-CSP the goal is simply to satisfy as may constraints as possible, additionally the domains are the variables are numerical (e.g. the real numbers).

While this problem has certainly been considered before, we were only able to find one previous work that directly address this issue.

In this paper the authors propose a solution to numerical MAX-CSP. In particular where $D_i = \mathbb{R}$ and each constraint $C_j = a_j^T x \leq b_j, a_j \in \mathbb{R}^n, b_j \in \mathbb{R}$. That is, given a set of linear inequalities, satisfy as many as possible. In the paper the authors propose a specific algorithm for this problem based on branch and bound techniques. Unfortunately, this algorithm is worst case exponential time. In fact, this problem is easy to formulate as a mixed integer program,

$$\begin{aligned} \min_{x,z} \quad & 1^T z \\ \text{s.t.} \quad & Ax - \epsilon z \leq 0 \\ & w \in \mathbb{R}^n \\ & z \in \{0, 1\}^m \end{aligned}$$

Where ϵ is some large number. The general problem of mixed integer linear programming is of course NP-Hard in general (by reduction to 0-1 integer programming which is NP-Complete). Hence this numerical MAX-CSP is likely not to admit an efficient solution.

3.0.2 Gradient Descent

Gradient descent is an optimization technique which has gained much attention in recent years due to the stochastic variant being used in training neural networks for machine learning tasks such as computer vision and speech recognition. Despite this focus on machine learning, gradient descent is a general purpose optimization procedure, capable of optimizing arbitrary differentiable functions. Many variations have been proposed in recent years,

however each technique uses the same basic idea. Given a function f and current point x , compute $f(x)$. Next compute the gradient $\nabla(f(x))$ and update x taking a step in a direction of $-\eta\nabla(f(x))$. Repeat this process until the minima of the function is found.

If $f(x)$ is convex, then gradient descent is guaranteed to find an optimal solution. If $f(x)$ is not convex, the gradient descent will converge to a local minimum. Although in the non-convex case, gradient descent is not guaranteed, and frequently won't, find the global minimum, it is useful for finding minima of functions which are cannot be handled by typical solvers. This makes it an invaluable tool for optimizing functions which don't admit exact optimization techniques such as those used in linear programs and quadratic programs.

4 Motivation

Numerical MAX-CSPs have a wide variety to applications in such as debugging infeasible linear programs, however we are interested in one particular problem, which is that of tuning the weights for full search engines.

Full text search engines are used in a plethora of applications. These search engines (such as Apache Lucene) are built for efficient retrieval of top-k documents based on a TF/IDF based scoring metric which is dot product between two sparse vectors $q^T d = \text{score}[[1]][1][[2]][2][[3]][3]$. While this default scoring gives decent results out of the box, it is frequently augmented by re-weighting the query q , with some weight vector w changing to scoring to be $(q \odot w)^T d = \text{score}$. This allows for boosting of certain terms or fields to improve the quality of search results while not changing the underlying search algorithm.

The problem we will address in this project is finding a good weight vector w . In particular our problem setting is as follows. We are giving a set of query vectors $Q = \{q_1, \dots, q_n\}$. For each of these query vectors q_i we are given a set of k retrieved document vectors with labels $R_i = \{(d_{i,1}, l_{i,1}), \dots, (d_{i,k}, l_{i,k})\}$, $l \in \{\text{relevant}, \text{irrelevant}\}$. We wish to find a weight vector w that minimizes the number of irrelevant documents that are examined before the relevant result is found. This problem can be formulated as a mixed integer linear program (MILP) of the following form

$$\begin{aligned} \min_{w,z} \quad & 1^T z \\ \text{s.t.} \quad & Aw - \epsilon z \leq \gamma \\ & w \geq 0 \\ & z \in \{0, 1\} \end{aligned}$$

5 Baseline Solutions

As mentioned above we can formulate our problem as a mixed integer linear program with binary constraints. In addition, to the mixed integer formulation, we also create an LP relaxation of the problem as an approximation of the MILP, namely,

$$\begin{aligned} \min_{w,z} \quad & 1^T z \\ \text{s.t.} \quad & Aw - \epsilon z \leq \gamma \\ & w \geq 0 \\ & 0 \leq z \leq 1 \end{aligned}$$

The rationale behind this is that the LP relaxation can be solved efficiently and provide a deterministic approximation of the MILP. Although the solution could be arbitrarily bad, we find that the LP typically finds a reasonable solution in much less time relative to the MILP.

6 Problem with Baseline Solutions

While modeling the MAX-CSP problem as a MILP is simple solution (at least programming-wise), it has a few major issues, namely that depending on the input, the solver time could be exponential in the number of constraints. Moreover, even commercial solvers can struggle with large problem sizes, either taking large amounts of time to

produce any feasible solution or having prohibitive memory requirements. For our motivating example, we would ideally feed in many labeled queries, each producing tens of constraints, hence exponential runtime in the number of constraints greatly reduces the usefulness of the solution.

7 Our Solution

To remedy problem of runtime we propose a new solution based on gradient descent. The high level idea of the solution is simple, begin with a random weight feasible vector w and then perform gradient descent until the local minima is found. More specifically we have,

$$f(x) = \text{HardSigmoid}(Aw - \gamma)$$

Where,

$$x \in R, \quad \text{HardSigmoid}(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ 1 & \text{if } x \geq 1 \\ x & \text{otherwise} \end{cases}$$

We

8 Experiments

We ran our experiments on a server with an AMD Ryzen Threadripper 2950X (16c/32t) processor with 64GB of RAM. Our MAX-CSP instances are generated from real world datasets, with the number of constraints ranging from **TODO : max num const** to **TODO : max num const**.

9 Discussion

10 Conclusions